

Express Mailing Label No. ET580099288US

PATENT APPLICATION
Docket No. DE920000124US1
File: 1200.2.44

UNITED STATES PATENT APPLICATION

of

Christian Bolik

Peter Gemsjaeger

and

Klaus Schroiff

for

**Method and System for Scalable, High Performance
Hierarchical Storage Management**

BRIAN C. KUNZLER
ATTORNEY AT LAW
10 WEST 100 SOUTH SUITE 425
SALT LAKE CITY, UTAH 84101

ET 580099288US

1 **Method and System for Scalable, High Performance**
2 **Hierarchical Storage Management**

3
4 **BACKGROUND OF THE INVENTION**

5 **1. The Field of the Invention**

6 The invention generally relates to hierarchical storage management systems and
7 more specifically to a method and system for managing an hierarchical storage management
8 (HSM) environment including at least one HSM server and at least one file server having
9 stored a managed file system, wherein the at least one HSM server and the at least one file
10 server are interconnected via a network and wherein digital data files are migrated
11 temporarily from the at least one file server to the at least one HSM server.
12

13 **2. The Relevant Art**

14 Hierarchical Storage Management (HSM) is used for freeing up more expensive
15 storage devices, typically magnetic disks, that are limited in size by migrating data files
16 meeting certain criteria, such as the age of the file or the file size, to lower-cost storage
17 media, such as tape, thus providing a virtually infinite storage space. To provide transparent
18 access to all data files, regardless of their physical location, a small "stub" file replaces the
19 migrated file in the managed file system. To the user this stub file is indistinguishable from
20 the original, fully resident file, but to the HSM system the stub file provides important
21 information such as where the actual data is located on the server.
22

23 An important difference between the views of a migrated file from the user's and the
24 HSM system's perspective is that the user doesn't see the new "physical" size of the file,
25 which after a file has been migrated is actually the size of the stub file, but still sees the
26 "logical" size, which is the same as the size of the file before it was migrated.

1 One implementation category of an HSM system makes use of a client/server setup,
2 where the client runs on the machine on which file systems are to be managed, and where
3 the server provides management of migrated data files and the included information.

4
5 Traditionally, an HSM system needs to perform the following tasks:

- 6 a) Determine which data files in the file system are eligible for migration
7 (referenced as "candidates"). In order to determine the "best" candidates (with
8 respect to their age and size), a full file system traversal is required;
- 9 b) Determine which previously migrated files have been modified in or removed
10 from the client file system so their migrated copies can be removed from the server
11 storage pool to reuse the space they occupied (referenced as "reconciliation"). To
12 accomplish this, usually a full file system tree traversal is necessary.

13
14 In case of insufficient available space in the client file system, data files need to be
15 migrated off the disk quickly to minimize application latency, herein referenced as
16 "automigration". If a managed file system runs out of space, all applications performing write
17 requests into this file system are blocked until enough space has been made available by
18 migrating files off the disk to satisfy their write requests. In traditional HSM systems, data
19 files in a managed file system are migrated serially, one file at a time.

20 An according data migration facility is disclosed in IBM Technical Disclosure
21 Bulletin, published June 1973, pp. 205 – 208. A supervisory controller is described for
22 automatic administration and control of a computer system's secondary storage resources.
23 A migration monitor is run-time event driven and acts as a first level event processor. The
24 migration monitor records events and summarizes a data migration activity. A migration task
25 is initiated by the migration monitor when a request is received. The migration task scans
26

1 through an inventory of authorized data on the system and invokes a given algorithm to make
2 the decision as to what data to migrate.

3 With the amount of data and the number of data files in a typical managed file system
4 increasing logarithmically over time as illustrated in Fig. 1, scalability of the HSM system
5 becomes an issue. Typical file system environments with such a behavior are those of
6 Internet providers handling the files of much more than thousands of customers, video
7 processing scenarios like those provided on a video-on-demand server, or weather forecast
8 picture processing where millions of high-resolution pictures are generated on a per day basis
9 by weather satellites. In those environments the number of files to be handled often exceeds
10 1 million and is continuously increasing.

11 For the above reasons, there exists a strong need to provide HSM systems which are
12 able to handle those very large file systems.

13 Most of the known HSM approaches traverse the complete file system in order to
14 gather eligible candidates for the automigration to remote storage. This system worked well
15 in rather small environments but are no longer usable for current file system layouts due to
16 the excessive processing time for millions of files. Therefore it is required to provide a more
17 scalable mechanism consuming less system resources.

18 A known HSM approach addressing an above migration scenario and disclosed in
19 U.S. Patent No. 5,832,522 proposes a placeholder entry (stub file) used to retrieve the status
20 of a migrated data file. In particular, a pointer is provided by which a requesting processor
21 can efficiently localize and retrieve a requested data file. Further, the placeholder entry
22 allows to indicate migration of a data file to a HSM server.

23 Another approach, a network file migration system is disclosed in U.S. Patent No.
24 5,367,698. The disclosed system comprises a number of client devices interconnected by a
25 network. A local data file storage element is provided for locally storing and providing
26 access to digital data files stored in one or more of the client file systems. A migration file

1 server includes a migration storage element that stores data portions of files from the client
2 devices, a storage level detection element that detects a storage utilization level in the storage
3 element, and a level-responsive transfer element that selectively transfers data portions of
4 files from the client device to the storage element.

5 Known HSM applications traverse the complete file system tree in order to gather
6 eligible candidates for the automigration to a remote storage. This system worked well in
7 rather small environments but are no longer usable for current file system layouts due to the
8 excessive processing time for millions of files. A complete tree traversal disadvantageously
9 impedes scalability both in terms of duration and resource requirements, as both numbers
10 grow logarithmically with the number of files in a file system. Furthermore, serial
11 automigration is often not capable of freeing up space quickly enough to satisfy today's
12 requirements. Therefore it is required to provide a more scalable mechanism consuming less
13 system resources.

14 Due to the ever increasing size of storage volume as well as the pure number of
15 storage objects makes it more and more difficult for a data management application to
16 provide its service without an increasing need for more system resources which is obviously
17 not desirable.

OBJECT AND BRIEF SUMMARY OF THE INVENTION

The hierarchical storage management system of the present invention has been developed in response to the present state of the art, and in particular, in response to the problems and needs in the art that have not yet been fully solved by currently available hierarchical storage management systems. Accordingly, it is an overall object of the present invention to provide a hierarchical storage management system that overcomes many or all of the above-discussed shortcomings in the art.

The underlying concept of the invention is, instead of attempting to find the "best" migration candidates all at once, to scan the file system only until a certain amount of migration candidates have been found. Further, the idea is that the process for determining candidates waits until one of two events to happen, namely until a specified wait interval expires or until an automigration process starts. The candidate determination process advantageously can resume the file system scan at the point where it stopped a previous scan and continue to look for migration candidates, again until a certain amount of candidates has been found.

The particular step of scanning the managed file system only until having detected a prespecified amount of migration candidate files advantageously enables that migration candidates are made available sooner to the migration process wherein migration can be performed as an automigration process not requiring any operator or user interaction. As the at least one attribute the file size and/or a time stamp of the file can be used.

In one embodiment, the automigration process is performed by a master/slave concept where the master controls the automigration process and selects at least one slave to migrate candidate data files provided by the master.

Another embodiment comprises the additional steps of ranking and sorting the candidate data files contained in the at least one list for identifying candidate data files, in particular with respect to the file size and/or time stamp of the data files contained in the at

1 least one list for identifying candidate data files. Hereby the order of candidate data files to
2 be migrated can be determined.

3 In particular, the proposed mechanism therefore makes the candidates determination
4 process practically independent from the number of files in the file system and from the size
5 of the file system. The invention therefore allows parallel processing of determination of
6 candidate data files for the migration and the automigration process itself.

7 In addition, the automigration process generates a unique identifier to be stored on
8 the HSM server that allows a direct access to migrated data files during a later reconciliation
9 process.

10 The proposed scanning process therefore significantly reduces resource requirements
11 since e.g. the storage resources for the candidate file list and the required processing
12 resources for managing the candidate file list are significantly reduced. In addition, the
13 scanning time is also reduced significantly.

14 The basic principal behind this invention is dropping the requirement of 100%
15 accuracy for the determination of eligible migration candidates. Rather than looking for an
16 analysis based on a complete list of migration candidates we can assume that the service is
17 also functional based on a certain subset of files within a managed file system.

18 Thereupon, the invention allows for handshaking between the process for determining
19 or searching migration candidates and the process of automigration.

20 As a result, the invention provides scalability and significant performance
21 improvement of such an HSM system. Thereupon secure synchronization or reconciliation
22 of the client and server storage without need of traversing a complete client file system is
23 enabled due to the unique identifier.

24 According to an embodiment, at least two lists for identifying candidate data files are
25 provided, whereby the first list is generated and/or updated by the scanning process and
26 whereby the second list is used by the automigration process. The automigration process

1 gathers the first list from the scanning process when all candidate data files of the second list
2 are migrated. Both lists are worked on in parallel thus revealing parallelism between
3 scanning and automigrating.

4 It is further noted, that besides the above described 'migrated' state, also a
5 'premigrated' state for data files in the managed file system can be used for which the
6 migrated copy stored on the HSM server is identical to the resident copy of the data file in
7 the managed file system.

8 These and other objects, features, and advantages of the present invention will
9 become more fully apparent from the following description and appended claims, or may be
10 learned by the practice of the invention as set forth hereinafter.

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

BRIEF DESCRIPTION OF THE DRAWINGS

In order that the manner in which the advantages and objects of the invention are obtained will be readily understood, a more particular description of the invention briefly described above will be rendered by reference to specific embodiments thereof which are illustrated in the appended drawings. Understanding that these drawings depict only typical embodiments of the invention and are not therefore to be considered to be limiting of its scope, the invention will be described and explained with additional specificity and detail through the use of the accompanying drawings in which:

Fig. 1 is a block diagram showing a typical hierarchical storage management (HSM) environment where the present invention can be applied to;

Fig. 2 illustrates the known logarithmic increase of the amount of data and the number of data files in a typical managed file system;

Fig. 3 is flow diagram for illustrating the basic mechanism of managing an HSM system according to the invention;

Fig. 4 is another flow diagram for illustrating the basic mechanism of reconciling a managed file system migrated from a file server to an HSM system;

Fig. 5 is another flow diagram showing a base logic of an automigration environment in accordance with the invention; and

Figs.6a,b illustrate a preferred embodiment of the mechanism according to the invention.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

Fig. 1 shows a typical file server 101 that manages one or more file systems 102. Each file system is usually organized in more or less complicated and more or less deeply nested file trees 103. The file server 101 is connected via a network 104, usually a Local Area Network (LAN) or a Wide Area Network (WAN), to another server machine 105 that contains an HSM server 106. The server machine 105 has one or more external storage devices 107, in this example tape storages, attached to. The HSM server 106 stores data, migrated from the file server 101 to the tape storages 107.

Fig. 2 illustrates that the amount of data and the number of data files in a typical managed file system is increasing logarithmically and is discussed beforehand.

The flow diagram depicted in Fig. 3 illustrates the basic mechanism of managing an HSM system according to the invention. In step 200, an amount of files, e.g. the number of files or the entire size of multiple files, for which a scan in the file system shall be performed, is pre-specified. Based on that pre-specified amount, at least part of the file system is scanned 201. It is an important aspect of the invention that not a whole file system is scanned through but only a part of it determined by the pre-specified amount.

In a next step 202, based on one or more attributes like the file size or a time stamp for the file (file age or the like), candidate files to be migrated from the file server to the HSM server are determined. The determined candidate files are put into a list of candidates 203. It is noteworthy hereby that, in another embodiment of the invention, two lists are provided. Such an embodiment is described hereinbelow in more detail.

Step 204 is an optional step (indicated by the dotted line) where the data files contained in the candidate list are additionally ranked in order to enable that the following selected files to be migrated can be migrated in a particular order.

In parallel to the steps 200 – 204 described above, the file system is monitored 205 and the current status of the file system is determined 206. In step 207, an automigration of

1 selected and allegedly ranked candidate data files is initiated or triggered by the determined
2 file system status. For the details of that file system status it is referred to the following
3 description.

4 After the automigration has been initiated, it is performed 208 by physically
5 transferring data files to the HSM server and, in particular, a unique identifier is assigned to
6 each migrated file. The concept and meaning of that unique identifier (ID) will become more
7 evident from the following parts of the description. Finally the unique identifier is sent to the
8 HSM server.

9 Now referring to the flow diagram depicted in Fig. 4, the basic mechanism of
10 reconciling a managed file system migrated from a file server to an HSM system, in
11 accordance with the invention, shall be illustrated. In a first step 301, a list of already
12 migrated data files is transferred via the network from the HSM server. The transferred list,
13 in particular, includes the unique identifier generated in the process described referring to
14 Fig. 3. Then a reconciliation process queries 302 the transferred list of migrated files and
15 compares 303 the migrated files, which are identified by their corresponding unique
16 identifier (ID) with the corresponding files contained in the managed file system. Finally, the
17 reconciliation process accordingly updates 304 the managed data on the HSM server.

18 The flow diagram depicted in Fig. 5 shows a base logic of an automated HSM
19 environment. A monitor daemon 501 starts a master scout process 502 and continuously
20 monitors one or more file system. The master scout process 502 starts one slave scout
21 process 503 per file system. Each slave scout process 503 scans its file system for candidate
22 data files to be migrated.

23 If the monitor daemon 501 detects that the file system has exceeded its threshold
24 limits, it starts a master automigration process 504, described in more detail hereinbelow.
25
26

1 If the value for a reconcile interval has exceeded, a reconciliation process 505 is started by
2 the monitor daemon 501. The reconciliation process 505 is also described in more detail in
3 the following.

4 The flow diagrams depicted in Fig. 6a and 6b illustrate a preferred implementation
5 based on independent migration candidates pools 601, 602 for the automigration 603 and
6 scanning process 604, the latter often (and in the following) referred to as "scout" process.

7 In this embodiment, the automigrator 603 is activated by another process - e.g. a
8 monitor process that tracks file system events and takes appropriate measurements if certain
9 thresholds are exceeded. The automigration 603 then starts to migrate 605 migration
10 candidates to a remote storage as long as some defined threshold is exceeded. Prior to
11 migrating 605 the files, the automigration process 603 performs management class (MC)
12 checks 606 with the HSM server to find out whether a potential migration does not violate
13 HSM server side rules.

14 If the automigration process 603 runs out of candidates, i.e. the list of identified
15 candidates 602 is used up, it sets 607 a flag to signal a request to the scout process 604 in
16 order to obtain a new list 601 of candidates. The scout process 604 receives 608 the flag and
17 moves 609 the newly generated list 601 to the automigrator 603, setting 609 another flag to
18 signal the automigrator 603 to continue with migrating files.

19 The scout process 604 itself starts to collect 610 new migration candidates. After
20 completion of the scanning, the scout process 604 will wait until it receives another signal
21 by the automigrator or by exceeding a definable value CANDIDATESINTERVAL 611. The
22 value CANDIDATESINTERVAL 611 defines the time period during the scout process 604
23 remains sleeping in the background after an activity phase.

24 In the latter case of exceeding the CANDIDATESINTERVAL 611, it starts
25 optimizing its candidates list with another scan. I.e. in case of not receiving a signal from
26 the automigration process, in order to improve the quality of the candidates list scout process

1 starts at each CANDIDATESINTERVAL 611 to scan for a new bunch of candidates. That
2 bunch of candidates is defined by another value MAXCANDIDATES 612 that defines a
3 number of required candidates following candidates criteria. Combined with the existing
4 migration candidates list 601 the scout process 604 can either collect all candidates or just
5 take the "best" subset in order to limit the required storage space. Thus the scout process
6 traverses the managed file system in order to find eligible candidates for automigration.
7 Rather than traversing the complete file system it stops as soon as MAXCANDIDATES 612
8 eligible candidates were found. Hereafter the process either waits for a dedicated event from
9 the automigration process or sleeps till CANDIDATESINTERVAL 611 time has passed.

10 The above scout process has the following advantages:

- 11 · Minimal consumption of system resources (memory, processing time) required
12 to find eligible candidates;
- 13 · highly scalable with minimal dependencies regarding the number of objects
14 within a file system;
- 15 · increasing candidates quality in times of normal file system activity.

16
17 As a possible disadvantage, it is possible that the potentially best migration
18 candidates based on the selection strategy are not used by the automigration process because
19 the scout process has not yet traversed the corresponding subtree. Nevertheless, the above
20 advantages considerably exceed the disadvantages.

21 In the following, the different process steps of the whole migration mechanism
22 proposed by the invention is described in more detail.
23
24
25
26

1 Candidates Determination

2 Avoiding full file system traversals:

3 Instead of attempting to find the "best" migration candidates in one shot, the file
4 system is scanned only until a certain number of migration candidates have been found.
5 Then, the candidates determination process waits for one of two events to happen:

- 6 · a specified wait interval expires, or
7 · automigration starts.

8 In this event, the process resumes the file system scan at the point where it left off and
9 continues to look for migration candidates, again until a certain number of candidates has
10 been found. These candidates are merged into the existing list of candidates and then
11 "ranked" for quality (with respect to age and size), thus incrementally improving the quality
12 of migration candidates in the system.

13 The benefit of this approach is that migration candidates are made available sooner
14 to the automigration process, and significantly reduced resource requirements, making the
15 candidates determination process practically independent from the number of files in the file
16 system and from the size of the file system.

17
18 Quick eligibility check:

19 A file can be eligible for migration only if it is not yet migrated. On file systems that
20 don't provide an XD SM API (X/Open Data Storage Management API), such as AIX JFS, the
21 migration state typically needs to be determined by reading a stub file. In order to limit the
22 number of files that the candidates determination process needs to read, usually only those
23 files are read whose physical size meets the criteria for being a stub file, but even then the
24 performance impact on file systems with a high percentage of migrated files is significant
25 as the read/write head of the hard disk constantly needs to jump back and forth between the
26 inode area of the file system and the actual data blocks. To address this, the present invention

1 proposes to require all stub files to have a certain characteristic, such as a specific physical
2 file size. The candidates determination process, then, can assume that all files whose
3 physical size matches the stub file size are migrated and disregard them from further
4 eligibility checking that would require reading the stub file. This will exclude resident files
5 whose size make them appear like stub files from migration, but the assumption is that the
6 percentage of such files in a typical file system is small enough to make this a viable
7 simplification.

8 In addition, the automigration process signals the need for additional migration
9 candidates. Once the file system exceeds a certain fill rate or runs out-of storage capacity the
10 automigration process gets started - usually initiated by the supervising daemon running
11 permanently in the background.

12 Hereby it consumes migration candidates from a dedicated automigration pool and signals
13 the scout process to dump his set of migration candidates to disk or to transfer it into a
14 migration queue via shared memory. Based on the newly dumped candidates list the
15 automigration process can now start to migrate data to the remote HSM server - preferably
16 multithreaded and via multiple processes where each migrator instance cares about a certain
17 set of files.

18 In order to guarantee maximum concurrency, the scout process can immediately start
19 to scan for new migration candidates after transferring his current list to the automigration
20 process. The immediate generation of a new candidates list insures that the automigration
21 process does not run out-of migration candidates or minimizes the wait time. Under normal
22 conditions new candidates are found much faster than the network transfer of the already
23 found candidates so we can assume that this is no bottleneck in this environment.
24
25
26

1 Automigration

2 Parallel Automigration

3 To lift the scalability limitations of the traditional serial automigration, the present
4 invention proposes a master/slave concept to facilitate parallel automigration of files in the
5 same file system. In this concept, a master automigration process reads from a list of
6 migration candidates created by the candidates determination process and dispatches entries
7 from this file to a certain number of automigration slaves ("migrators"). These slaves migrate
8 the file they are assigned to the HSM server, and then are available again for migrations as
9 assigned by the master process.

10 The essential benefit is the scalability of the speed by which files can be migrated off
11 the file system, by defining the number of parallel working automigration slaves. The
12 complete control of the automigration process remains sequential (master automigration
13 process), so that no additional synchronization effort is required, as it would be like in other
14 typical parallel working systems. The "real work", the migration of the files itself, that
15 consumes most of the time during the whole automigration process, is parallelized.

16
17 Reconciliation

18 Immediate Synchronization:

19 To reconcile a client/server HSM system, the HSM client, according to the prior art,
20 has to perform the following steps:

21
22 Retrieve the list of migrated files for a given file system from the HSM server
23 (the "server list") and

24
25 Traverse the file system tree, marking each unmodified migrated file as
26 "found" in the server list.

1 When tree traversal is completed, all files in the server list not marked "found" will
2 be marked for removal from a server storage pool, as they were either removed from the
3 client file system, or their client copy was modified, thus invalidating the server copy. The
4 reconciliation processing known in the prior art therefore requires a full file system tree
5 traversal, which poses the scalability problems described above. To avoid the need for a full
6 traversal, the invention proposes the following processing:

7 • When migrating files, the HSM client stores a unique, file system-specific
8 identifier (the "file ID") with the file on the HSM server;

9 • during reconciliation, the HSM client retrieves the list of migrated files, in
10 particular by use of the unique ID stored in the list or array, from the server as before,
11 but now the server list includes the file id for each entry;

12 • for each entry from the server list received, the HSM client invokes a
13 platform-specific function that returns the file attributes of a file identified by its file
14 id. On IBM AIX (UNIX derivate) this makes use of the vfs_vget VFS entry point,
15 which should be invoked so that it reads the attributes directly from the underlying
16 physical file system to avoid having to read the stub file, whereas on DMAPI- enabled
17 file systems the dm_get_fileattr API is used;

18 • if the attributes could be determined and match with those stored in the server
19 list, processing continues with step 3 until all entries have been received. Otherwise
20 the entry will be added to a list in client memory that will be used to mark files for
21 removal on the server (the "remove list");

22 • when all entries from the server list have been received and processed, the
23 HSM client loops through the remove list, and marks each of them for removal from
24 the server storage pool.

1 Quick premigration check:

2 In addition to the "migrated" and "resident" states of a file, some HSM systems provide a
3 third state: "premigrated". A file is "premigrated" when its copy on the server (after
4 migration) is identical to the (resident) copy of the file in the client file system. This is the
5 case for instance immediately after a migrated file is copied back to the local disk: the file
6 is resident, but its migrated copy is still present in the server storage pool, and both copies
7 are identical.

8 The benefit of the premigration state is that such files can be migrated simply by
9 replacing them with a stub file, without having to migrate the actual data to the HSM server.
10 On file systems that don't provide the XDSM API the HSM client needs to keep track of the
11 premigrated files in a look-aside database (referenced as "premigration database"), as
12 premigrated files don't have an associated stub file that could be used to store premigration
13 information.

14 Those HSM clients, that rely on a look-aside database, need to traverse the local file
15 system to verify the contents of the premigration database. However, making use of the
16 same principle proposed in the previous section "Immediate Synchronization", the need for
17 a full tree traversal can be removed here as well by storing a unique file id for each
18 premigrated file in the premigration database, and then perform a direct mapping from its
19 entries into the file system. Entries whose mapping is no longer successful can be removed
20 from the premigration database.

21 Finally it is emphasized that combined with one another, the proposed measures
22 resolve the most pressing scalability problems and performance bottlenecks present in
23 traditional client/server-based HSM systems.

24 What is claimed and desired to be secured by United States Letters Patent is:
25
26